

1 SYSTEM AND METHOD FOR RECOVERING FROM MEMORY FAILURES IN
2 COMPUTER SYSTEMS

3
4 BACKGROUND OF THE INVENTION

5 1. Field of the Invention

6 The present invention relates generally to systems and methods for recoverable
7 programming, and more particularly to a recoverable programming system and method
8 for memory system failures in multi-processor computer systems.

9 2. Discussion of Background Art

10 Demand for increased performance and high availability of commodity computers
11 is increasing with the ubiquitous use of computers and the Internet services which serve
12 them. While commodity systems are tackling the performance issues, availability has
13 received less attention. It is a common belief that software (SW) errors and administration
14 time are, and will continue to be, the most probable cause of the loss of availability.
15 While such failures are clearly commonplace, especially in desktop environments, it is
16 believed that certain other hardware (HW) errors are also becoming more probable.

17 Processors, caches, and memories are becoming larger, faster and more dense,
18 while being increasingly used in ubiquitous and adverse environments such as at high
19 altitudes, in space, and in industrial applications. Articles, such as Ziegler, J. F., et al.,
20 "IBM Experiments in Soft Fails in Computer Electronics (1978-1994)", IBM Journal of
21 R&D, vol 40, no 1, pp 3-18, January 1996, and Ziegler, J. F., "Terrestrial Cosmic Rays",
22 IBM Journal of R&D, vol 40, no 1, pp 19-40, January 1996, have shown that these
23 changes will lead to increased transient errors in CMOS memory due to the effects of

1 cosmic rays, approximately 6000 FIT (1 FIT equals 1 failure in 10^9 h) for one 4Mbit
2 DRAM.

3 Tandem (see, Compaq Corporation, "Data Integrity for Compaq NonStop
4 Himalaya Servers", White Paper, 1999) indicates that such errors also apply to processor
5 cores or on-chip caches at modern die sizes/voltage levels. They claim that processors,
6 cache, and main memory are all susceptible to high transient error rates. A typical
7 processor's silicon can have a soft-error rate of 4000 FIT, of which approximately 50%
8 will affect processor logic and 50% the large on-chip cache. Due to increasing speeds,
9 denser technology, and lower voltages, such errors are likely to become more probable
10 than other single hardware component failures. With the increasing evolution to larger
11 tightly interconnected commodity machines (such as Sun's Enterprise 10000 machines),
12 the probability of soft-errors and error containment problems increases further. Soft-error
13 probability increases not only due to increased system scale, but also due to an increased
14 number of components on the memory access path. Since the machines are tightly
15 coupled, memory path soft-errors introduce error containment problems which without
16 some form of soft-ware error containment can lead to complete loss of availability.

17 Techniques such as Error Correction Codes (ECC) and ChipKill (see, Dell, T. J.,
18 "A White Paper on the benefits of Chipkill Correct ECC for PC Server Main Memory",
19 IBM Microelectronics Division, Nov. 1997) have been used in main memories and
20 interconnects to correct some of these errors (90% for ECC). Unfortunately such
21 techniques, only help reduce visible error rates for semiconductor elements that can be
22 covered by such codes (large storage elements). With raw error rates increasing with
23 technological progress and more complicated interconnected memory subsystems, ECC is
24 unable to address all the soft-error problems. For example, a 1Gb memory system based

1 on 64Mbit DRAMs still has a combined visible error rate of 3435 FIT when using Single
2 Error Correct Double Error Detect (SECCDED) ECC. This is equivalent to around 900
3 errors in 10000 machines in 3 years. Unfortunately, current commodity hardware and
4 software provide little to no support for recovery from errors not covered by ECC
5 whether detected or not. Such problems have been considered by mainframe technology
6 for years, but in the field of commodity hardware, it is currently not cost effective to
7 provide full redundancy/support in order to mask errors. Therefore, the burden falls to
8 commodity hardware and the software using it to attempt to handle these errors for the
9 highest availability.

10 Most contemporary commodity computer systems, while providing good
11 performance, pay little attention to availability issues resulting from such errors. For
12 example, the IA-32 architecture supports only ECC on main memory rather than across
13 the system, requiring system reboot on errors not covered by this ECC. Consequently,
14 commodity software such as the OS, middleware and applications have been unable to
15 deal with the problem. Future commodity processor architectures may provide support to
16 detect and notify the system of such probable errors. For instance, upcoming IA-64
17 processors, while not recoverable in the general case, do offer some support with certain
18 limitations.

19 Availability in computer systems is determined by hardware and software
20 reliability. Hardware reliability has traditionally existed only in proprietary servers, with
21 specialized redundantly configured hardware and critical software components, possibly
22 with support for processor pairs (see, Bartlett, J., "A Nonstop Kernel", Proceedings of the
23 Eighth Symposium on Operating Systems Principles, Asilomar, Ca, pp 22-29, Dec.
24 1981), e.g. IBM S/390 Parallel Sysplex (see, Nick, J.M., et al., "S/390 Cluster

1 Technology: Parallel Sysplex", *IBM Systems Journal*, vol 36, no 2., pp 172-201, 1997),
2 and Tandem NonStop Himalaya (see, Compaq, Product description for Tandem Nonstop
3 Kernel 3.0. Download Feb. 2000, <http://www.tandem.com>).
4 Sysplex supports hot swap execution, redundant shared disk with fault-aware system
5 software for error detection and fail-over restart. Tandem supports redundant fail-over
6 lock-stepped processors with a NonStop kernel and middleware, which provide improved
7 integrity through the software stack. These systems provide full automatic support to
8 mask the effects of data and resource loss. They rely on reliable memory and fail-over
9 rather than direct memory error recovery. Another approach is fault containment and
10 recovery with "node" granularity. In these systems, each node has its own kernel. When
11 one node fails, the others can recover and continue to provide services. Systems of this
12 type include the early cluster systems (see, Pfister, G., "In Search of Clusters", Prentice
13 Hall, 1998), and NUMA architectures, such as Hive (see, Chapin, J., et al., "Hive: Fault
14 Containment for Shared Memory Multiprocessors," *Proc. of the 15th SOSP*, Dec.1995,
15 pp 12-25, and Teodosiu, D., et al., "Hardware Fault Containment in Scalable
16 Shared Memory Multiprocessors," *Proc. of the 24th ISCA*, pp 73-84, June 1997).
17 Hardware faults are difficult to catch and repeat.
18 Software reliability has been more difficult to achieve in commodity software even with
19 extensive testing and quality assurance. Commodity software fault recovery has not
20 evolved very far. Most operating systems support some form of memory protection
21 between units of execution to detect and prevent wild read/writes. But most commodity
22 operating systems have not tackled problems of memory errors themselves or taken up
23 software reliability research in general. Examples include Windows 2000 and Linux.
24 They typically rely on failover solutions, such as Wolfpack by Microsoft. A lot of work

1 has been undertaken in the fault-tolerant community regarding the problems of reliability
 2 and its recovery in software (see, Brown, N.S. and Pradhan, D.K. "Processor and
 3 Memory-Based Checkpoint And Rollback Recovery", IEEE Computer, pp 22-31, Feb.
 4 1993; Gray, J., and Reuter, A., "Transaction Processing: Concepts and Techniques,"
 5 Morgan Kaufmann, 1993; and Kermarrec, AM., et al., "A Recoverable Distributed
 6 Shared Memory Integrating Coherence and Recoverability", *Proc. of the 25th FTCS*, pp
 7 289-298, June 1995).

8 These include techniques such as checkpointing and backward error recovery. A lot of
 9 this work has been conducted in the context of distributed systems rather than in single
 10 systems. There are also techniques for efficient recoverable software components, e.g.
 11 RIO file cache (see, Chen, P.M., et al., "The Rio File Cache: Surviving Operating
 12 System Crashes", *Proc. of the 7th ASPLOS*, pp 74-83, October 1996), and Recoverable
 13 Virtual Memory (RVM) (see, Satyanarayanan, et al. "Lightweight Recoverable Virtual
 14 Memory". *Proc. SOSP*, pp 146-160, Dec. 1993).

15 Rio takes an interesting software-based approach to fault containment aimed at a
 16 fault-tolerant file cache, but with general uses. By instrumenting access to shared data
 17 structures with memory protection operations, wild access to the shared data structures
 18 becomes improbable.

19 Other methods for handling memory errors include a try-except block solution. In
 20 general, the try-except mechanism itself is not sufficient to handle memory failures. The
 21 saved state needed for memory failures is more extensive (as an example, for IA-64
 22 architecture) than what can be obtained by try-except. Thus saving state is an expensive
 23 operation in terms of system overhead.

1

2

8

18

23

1 now be responded to and remedied without rebooting the computer. The present
2 invention will succeed in responding to memory errors much more effectively than
3 standard machine check abort handles.

4 The present invention is particularly applicable to O/S level code which can not
5 otherwise be restarted in response to memory errors without rebooting. When the present
6 invention is incorporated within application level code, the present invention also enables
7 the application to recover from memory errors, instead of otherwise being shut down and
8 restarted.

9 These and other aspects of the invention will be recognized by those skilled in the
10 art upon review of the detailed description, drawings, and claims set forth below.

11

1 BRIEF DESCRIPTION OF THE DRAWINGS

2 Figure 1 is a dataflow diagram of a system for recovering from memory access
3 failures;

4 Figure 2 is a flowchart of a first embodiment of a method for recovering from
5 memory access failures;

6 Figure 3 is a flowchart of a second embodiment of the method for recovering from
7 memory access failures; and

8 Figure 4 is a flowchart of a third embodiment of the method for recovering from
9 memory access failures.

10

1 DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENT

2 Figure 1 is a dataflow diagram of a system 100 for recovering from memory
3 access failures. The system 100 includes a memory 102, a memory controller 104, error
4 logging registers 106, and a central processing unit (CPU) 108, coupled together by a bus
5 110. The CPU 108 is controlled by software 112. The software 112 is preferably
6 included within the system's 100 operating system, however the software 112 could also
7 be instantiated within an application program as well. The software 112 configures the
8 memory controller 104 and has access to the error registers 106. Operation of the system
9 100 is discussed with respect to method Figures 2, 3, and 4.

10
11 Figure 2 is a flowchart of a first embodiment of a method 200 for recovering from
12 memory access failures. The method 200 begins with step 202 where the software 112
13 identifies a predetermined critical computer instruction sequence about to be executed by
14 the CPU 108, which includes a memory access request. The predetermined critical
15 computer instruction sequence can be part of a set of instruction sequences, identified by
16 the software 112 designer, for which error recovery is required. While the critical
17 computer instruction sequence discussed below include a memory access request, those
18 skilled in the art will know that concepts discussed with respect to recovery from a
19 memory access request error can be applied to other recovery critical instruction
20 sequences which would otherwise require rebooting of the system 100 to recover from.
21 Thus, the present invention is particularly applicable to O/S level code which can not
22 otherwise be restarted in response to memory errors without rebooting. When the present
23 invention is incorporated within application level code, the present invention enables the
24 application to recover from memory errors, instead of otherwise being shut down and

1 restarted. The present invention may also be used on non-critical computer instruction
2 sequences and for non-memory related errors.

3 In step 204, the software 112 then instructs the memory controller 104 to mask
4 any raised machine check abort (MCA) handle. Next in step 206, the CPU 108 executes
5 the memory access request. The memory controller 104 logs any memory access error in
6 the error logging register 106, in step 208. Next, in step 210, the software 112 polls the
7 error logging register 106 for any memory access errors, during execution of the
8 instruction sequence.

9 In step 212, the software 112 raises exceptions and updates pointers, if a memory
10 access error was logged during execution of the instruction sequence. The exceptions
11 perform various diagnostic functions in response to the memory error. The housekeeping
12 functions may include system recovery, memory management, and other reset procedures.
13 Pointers are updated when during memory error diagnosis, there are indications that a
14 portion or sector of the memory 102 may be physically damaged or corrupt.

15 Depending upon the memory access error which occurred, the software 112 may
16 command the CPU 108 to re-execute the memory access request, in step 214. The
17 software 112 will command the CPU 108 to re-execute the memory access request if the
18 memory access error detected is most likely due to a transitory error condition, which is
19 not likely to occur again. On the other hand, if the memory access error suggest that the
20 memory 102 itself is physically damaged, the software 112 will not instruct the CPU 108
21 to re-execute the memory access request. In step 216, the software 112 instructs the
22 memory controller 104 to enable the MCA handle.

23

1 Figure 3 is a flowchart of a second embodiment of the method for recovering from
2 memory access failures. The method 300 begins with step 302 where the software 112
3 identifies a predetermined critical computer instruction sequence about to be executed by
4 the CPU 108, which includes a memory access request.

5 In step 304, the software 112 then checkpoints a predetermined set of system data
6 necessary to recover should the memory access request fail. Checkpointing is component
7 of a transactional paradigm in which permanent modifications to system data are not
8 made until all associated operations within the transaction have been successfully
9 committed. Thus if during a transaction, such as the memory access request, an error
10 occurs, the system data stored during the checkpoint can be restored.

11 In step 306, the software 112 then instructs the memory controller 104 to mask
12 any raised machine check abort (MCA) handle. In step 308, the CPU 108 executes the
13 memory access request. The memory controller 104 logs any memory access error in the
14 error logging register 106, in step 310. Next, in step 312, the software 112 polls the error
15 logging register 106 for any memory access errors, during execution of the instruction
16 sequence. If a memory access error is logged during execution of the instruction
17 sequence, the software 112: raises exceptions and updates pointers, in step 314; recovers
18 the checkpointed system data, in step 316; and restores the system data, in step 318

19 As discussed above, with reference to Figure 2, depending upon the memory
20 access error which occurred, the software 112 may command the CPU 108 to re-execute
21 the memory access request, in step 320. In step 322, the software 112 instructs the
22 memory controller 104 to enable the MCA handle.

23

1 Figure 4 is a flowchart of a third embodiment of the method for recovering from
2 memory access failures. The method 400 begins with step 402 where the software 112
3 identifies a predetermined critical computer instruction sequence about to be executed by
4 the CPU 108, which includes a memory access request. In step 404, the software 112
5 then instructs the memory controller 104 to mask any raised machine check abort (MCA)
6 handle. In step 406, the CPU 108 executes the memory access request. The memory
7 controller 104 logs any memory access error in the error logging register 106, in step 408.

8 In step 410, the memory controller 104 sets data returned in response to the
9 memory access request equal to a set of predefined fake data, if a memory access error is
10 logged during execution of the instruction sequence. The software 112 has
11 preprogrammed the memory controller 104 to perform the functionality described in step
12 410. By setting the returned data to the predefined fake data in when a memory access
13 error occurs, corrupted data is not returned to the software, which might otherwise
14 necessitate a system reboot.

15 In step 412, the software 112 receives data returned in response to the memory
16 access request. In step 414, the method 400 skips to step 422, if the data returned in
17 response to the memory access request is not equivalent to the predefined fake data.
18 When the data returned is not equal to the fake data, the software 112 knows that no
19 memory access error has occurred, during execution of the instruction sequence, even
20 though the software 112 has not polled the error logging register. Thus, the polling step
21 can be eliminated, speeding up the memory access request.

22 In step 416, the software 112 polls the error logging register 106 for any memory
23 access errors, during execution of the instruction sequence. In step 418, the software 112

1 raises exceptions and updates pointers, if a memory access error was logged during
2 execution of the instruction sequence.

3 As discussed above, with reference to Figure 2, depending upon the memory
4 access error which occurred, the software 112 may command the CPU 108 to re-execute
5 the memory access request, in step 420. In step 422, the software 112 instructs the
6 memory controller 104 to enable any hardware raised MCA handles.

7
8 Another enhancement which may be applied to each of the three embodiments
9 discussed above, is to batch access to memory in large chunks whenever possible. By
10 batch accessing data, memory access errors are logged and polled for the entire batch.
11 This has implication on a granularity of the system 100 operation and is limited by
12 pointer manipulation.

13
14 While one or more embodiments of the present invention have been described,
15 those skilled in the art will recognize that various modifications may be made. Variations
16 upon and modifications to these embodiments are provided by the present invention,
17 which is limited only by the following claims.